

## TENSOR DECOMPOSITION MODELS

Why do we want to factorize a matrix/tensor?

- ↳ Compression / reduce model size
- ↳ Discover structure in data (PCA)
- ↳ Infer missing entries (completion)

### (I) MATRIX LOW RANK FACTORIZATION

Theorem: Let  $A \in \mathbb{R}^{m \times m}$

$\text{rank}(A) \leq R \iff \text{There exists } B \in \mathbb{R}^{m \times R} \text{ and } C \in \mathbb{R}^{R \times m}$   
such that  $A = BC$

↳ "rank  $R$  factorization of  $A$ "

$$\begin{matrix} m \\ \downarrow \\ A \end{matrix} = \begin{matrix} m \\ \downarrow \\ B \end{matrix} \begin{matrix} R \\ \downarrow \\ C \\ m \end{matrix}$$

$m^2$  parameters Vs.  $R(m+n)$  parameters

Sometimes, an approximate factorization is enough:

$$A \approx BC$$

$m \times m \quad m \times R \quad R \times n$

Low rank approximation problem: Given  $A \in \mathbb{R}^{m \times n}$ , and a target rank  $R$

$$\min_{\substack{B \in \mathbb{R}^{m \times R} \\ C \in \mathbb{R}^{R \times n}}} \|A - BC\|_F$$

Def (Frobenius norm):  $\|A\|_F = \sqrt{\sum_{i,j} (A_{ij})^2}$

Def (Singular Value Decomposition, SVD)

Any matrix  $A \in \mathbb{R}^{m \times m}$  ( $m \leq n$ ) can be decomposed into

① 
$$A = \underbrace{UDV^T}_{m \times m \quad m \times m \quad m \times m}$$

where  $U$  and  $V$  are orthonormal (i.e.  $U^T U = I$  and  $V^T V = I$ )  
and  $D$  is a diagonal matrix with non-negative entries.

The columns of  $U$  are the left singular vectors:  $u_1, u_2, \dots, u_m \in \mathbb{R}^m$

The columns of  $V$  are the right singular vectors:  $v_1, v_2, \dots, v_m \in \mathbb{R}^n$

The diagonal entries of  $D$  are the singular values:  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_m \geq 0$

Remark: We can rewrite the SVD (①):  $A = \sum_{i=1}^m \sigma_i u_i v_i^T$

If  $\text{rank}(A) = R$ , then  $\sigma_{R+1} = \sigma_{R+2} = \dots = \sigma_m = 0$  and  $A = \sum_{i=1}^R \sigma_i u_i v_i^T$

Theorem (Eckart-Young)

Let  $A \in \mathbb{R}^{m \times n}$  and let  $A = UDV^T = \sum_{i=1}^m \sigma_i u_i v_i^T$  be its SVD.  
Then the solution of

$$\min_{X \in \mathbb{R}^{m \times n}} \|A - X\|_F^2 \quad \text{subject to } \text{rank}(X) \leq R$$

is given by  $X^* = \sum_{i=1}^R \sigma_i u_i v_i^T$ . (truncated SVD)

Low rank approximation problem: Given  $A \in \mathbb{R}^{m \times n}$ , and a target rank  $R$

$$\min_{\substack{B \in \mathbb{R}^{m \times R} \\ C \in \mathbb{R}^{R \times n}}} \|A - BC\|_F$$

**SOLUTION**

Compute SVD of  $A = UDV^T$

keep the first  $R$  columns of  $U$  and  $V$  and the first  $R$  diagonal elements of  $D$ :

$$A = \begin{bmatrix} U_R & & \\ & \vdots & \\ & u_1 & u_2 & \dots & u_R & \dots & u_m \\ & & \vdots & & \vdots & & \vdots \\ & & & & & \ddots & \\ & & & & & & D_R \\ & & & & & & \sigma_1 & \sigma_2 & \dots & \sigma_R & \dots & \sigma_m \\ & & & & & & & & & & & & \sigma_m \\ & & & & & & & & & & & & \\ & & & & & & & & & & & & V_R^T \\ & & & & & & & & & & & & \begin{bmatrix} v_1^T \\ \vdots \\ v_R^T \\ \vdots \\ v_m^T \end{bmatrix} \end{bmatrix}$$

$$A \approx \underbrace{U_R}_{m \times R} \underbrace{D_R}_{R \times R} \underbrace{V_R^T}_{R \times m} = \underbrace{(U_R D_R)}_B \underbrace{(V_R^T)}_C$$

## II TENSOR NETWORKS (TN)

TN are graphs representing operations between tensors:

- Nodes represent tensors
- The arity (# of incoming edges / # of legs) of a node correspond to the order of the tensor:

$$v \in \mathbb{R}^d : \quad \begin{matrix} v \\ |_d \end{matrix} \quad A \in \mathbb{R}^{m \times n} : \quad \begin{matrix} -A- \\ m \quad n \end{matrix}$$

$$T \in \mathbb{R}^{d_1 \times d_2 \times d_3} : \quad \begin{matrix} T \\ d_1/d_2 \backslash d_3 \end{matrix}$$

- Edges represent contractions (summations):

+ Two matrices:  $\begin{matrix} -A- \\ m \quad n \end{matrix} \quad \begin{matrix} -B- \\ n \quad k \end{matrix}$

$$\left( \begin{matrix} -A- \\ m \end{matrix} \begin{matrix} B \\ k \end{matrix} \right)_{ij} = \sum_{k=1}^n A_{ik} B_{kj} = (AB)_{ij}$$

$$\hookrightarrow -A-B- = AB$$

+  $\underbrace{m \frac{d}{d} n} = \sum_{i=1}^d m_i n_i = \langle m, n \rangle$

+  $A \in \mathbb{R}^{m \times m}, \quad v \in \mathbb{R}^m \quad \left( \begin{matrix} -A- \\ m \end{matrix} v \right)_i = \sum_{j=1}^m A_{ij} v_j = (Av)_i$

+  $A \in \mathbb{R}^{m \times m} \quad \overbrace{\begin{matrix} A \\ m \end{matrix}} = \sum_{i=1}^m A_{ii} = \text{Tr}(A)$

+ Proof of  $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB) \quad (\neq \text{Tr}(BAC))$

$$\overbrace{\begin{matrix} A-B-C \end{matrix}} = \begin{pmatrix} A \\ B-C \end{pmatrix} = \begin{pmatrix} B \\ C-A \end{pmatrix}$$

+  $m \in \mathbb{R}^m, \quad v \in \mathbb{R}^m. \quad \left( \begin{matrix} m \\ m \end{matrix} \begin{matrix} v \\ m \end{matrix} \right)_{i,j} = m_i v_j = (m \circ v)_{i,j}$

$$\hookrightarrow \begin{matrix} m \\ | \end{matrix} \begin{matrix} v \\ | \end{matrix} = m \circ v$$

$$A \in \mathbb{R}^{m \times n}, \quad B \in \mathbb{R}^{n \times q} \quad \left( \begin{matrix} A \\ m \end{matrix} \begin{matrix} B \\ n \end{matrix} \begin{matrix} \\ q \end{matrix} \right)_{ijk\ell} = A_{ij} B_{k\ell}$$

$$+ \quad T \in \mathbb{R}^{d_1 \times d_2 \times d_3}, \quad A \in \mathbb{R}^{m \times d_2}, \quad B \in \mathbb{R}^{m \times d_3}, \quad C \in \mathbb{R}^{1 \times d_1}$$

$$\mathbb{R}^{d_1 \times m \times d_2} \ni \begin{array}{c} T \\ \diagdown d_1 \\ \diagup d_2 \\ A \\ | m \end{array} = T \times_2 A$$

$$\mathbb{R}^{d_1 \times d_2 \times m} \ni \begin{array}{c} T \\ \diagdown d_1 \\ \diagup d_2 \\ B \\ | m \end{array} = T \times_3 B$$

$$\mathbb{R}^{1 \times d_2 \times d_3} \ni \begin{array}{c} T \\ \diagdown d_1 \\ \diagup d_2 \\ C \\ | m \end{array} = T \times_1 C$$

$$+ \quad T \in \mathbb{R}^{d_1 \times d_2 \times d_3}, \quad A \in \mathbb{R}^{d_2 \times m}, \quad S \in \mathbb{R}^{m_1 \times m_2 \times m_3 \times m}$$

$$\begin{array}{c} d_1 \\ i \\ \diagdown d_2 \\ A \\ | m \\ \diagup m_1 \\ i_1 \\ \diagdown m_2 \\ i_2 \\ \diagup m_3 \\ i_3 \\ \diagdown m_4 \\ i_4 \\ \diagup m_5 \\ i_5 \end{array} = \sum_{j=1}^{d_2} \sum_{k=1}^m T_{i,j,i_2} A_{j,k} S_{i_3, i_4, i_5, k}$$

$$+ \quad m, n, w \in \mathbb{R}^d : \quad \mathbb{R} \ni \begin{array}{c} m \\ \diagdown d \\ n \\ \diagup d \\ w \end{array} = \sum_{i=1}^d m_i n_i w_i$$

### III 3 TENSOR DECOMPOSITION MODELS

#### 1) CP decomposition

$\hookrightarrow$  (CANDECOMP / PARAFAC)

$$\star \quad \begin{matrix} T \\ d_1 \diagdown d_2 \diagup d_3 \end{matrix} = \underbrace{\begin{matrix} R \\ A & B & C \\ d_1 \diagdown d_2 \diagup d_3 \end{matrix}}_{R} \quad \leftarrow \text{rank } R \text{ CP decomposition of } T$$

$A \in \mathbb{R}^{d_1 \times R}, B \in \mathbb{R}^{d_2 \times R}, C \in \mathbb{R}^{d_3 \times R}$

Def: The CP rank of a tensor  $T$  is the smallest  $R$  such that a rank  $R$  CP decomposition of  $T$  exists.

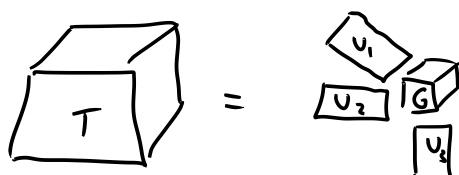
$$\star: T_{ijk} = \underbrace{\begin{matrix} T \\ A & B & C \\ i & j & k \end{matrix}}_{R} = \sum_{n=1}^R A_{in} B_{jn} C_{kn}$$

$$\text{If } A = \begin{pmatrix} a_1' & \dots & a_R' \\ 1 & \dots & 1 \end{pmatrix}, B = \begin{pmatrix} b_1' & \dots & b_R' \\ 1 & \dots & 1 \end{pmatrix}, C = \begin{pmatrix} c_1' & \dots & c_R' \\ 1 & \dots & 1 \end{pmatrix}$$

$$\text{then } T_{ijk} = \sum_{n=1}^R (a_n)_i (b_n)_j (c_n)_k = \sum_{n=1}^R (a_n \circ b_n \circ c_n)_{ijk}$$

$$\hookrightarrow T = \sum_{n=1}^R a_n \circ b_n \circ c_n$$

#### 2) TUCKER decomposition



$$T \in \mathbb{R}^{d_1 \times d_2 \times d_3}, \quad G \in \mathbb{R}^{R_1 \times R_2 \times R_3}, \quad U_i \in \mathbb{R}^{d_i \times R_i} \text{ for } i=1,2,3$$

$\hookrightarrow$  "core tensor",  $\hookrightarrow$  "factor matrices"

$$(Δ) \quad \begin{matrix} T \\ d_1 \diagdown d_2 \diagup d_3 \end{matrix} = \underbrace{\begin{matrix} G \\ R_1 \diagdown R_2 \diagup R_3 \\ U_1 & U_2 & U_3 \\ d_1 \diagdown d_2 \diagup d_3 \end{matrix}}_{R_1, R_2, R_3} \quad \leftarrow \text{rank } (R_1, R_2, R_3) \text{ Tucker decomposition of } T.$$

def: The multilinear rank (Tucker rank) of  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  is the smallest tuple  $(R_1, R_2, R_3)$  s.t. a rank  $(R_1, R_2, R_3)$  Tucker decomposition of  $T$  exists.

$\hookrightarrow$  The Tucker decomposition  $(Δ)$  can be written as  $T = G \times_1 U_1 \times_2 U_2 \times_3 U_3$

3) Tensor train decomposition (TT)

$T \in \mathbb{R}^{d_1 \times d_2 \times d_3 \times d_4}$ ,  $G^1 \in \mathbb{R}^{d_1 \times R_1}$ ,  $G^2 \in \mathbb{R}^{R_1 \times d_2 \times R_2}$ ,  $G^3 \in \mathbb{R}^{R_2 \times d_3 \times R_3}$ ,  $G^4 \in \mathbb{R}^{R_3 \times d_4}$

*"core tensors"*

$$\begin{array}{c} T \\ \diagdown d_1 \quad \diagup d_4 \\ \diagup d_2 \quad \diagdown d_3 \end{array} = G^1 \frac{R_1}{d_1} G^2 \frac{R_2}{d_2} G^3 \frac{R_3}{d_3} G^4$$

↳ rank  $(R_1, R_2, R_3)$  TT decomposition of  $T$ .

Def: The TT rank of a tensor  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3 \times d_4}$  is the smallest  $(R_1, R_2, R_3)$  such that rank  $(R_1, R_2, R_3)$  TT decomposition of  $T$  exists.

#### (IV) TENSOR LOW RANK APPROXIMATION

For  $T \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$  ( $d_1, d_2, \dots, d_N$  parameters)

	rank	# parameters	EXACT DECOMPOSITION	LOW RANK APPROXIMATION
CP	$R$	$R(d_1 + d_2 + \dots + d_N)$	NP-hard	NP
TUCKER	$(R_1, R_2, \dots, R_N)$	$R_1 R_2 \dots R_N + \sum_{i=1}^N d_i R_i$	P (easy)	NP
TT	$(R_1, \dots, R_{N-1})$	$d_1 R_1 + R_1 d_2 R_2 + R_2 d_3 R_3 + \dots + R_{N-2} d_{N-1} R_{N-1} + R_{N-1} d_N$	P	NP

#### Low rank approximation problem:

1) CP Given  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  and a target rank  $R$ :

$$\min_{\substack{A \in \mathbb{R}^{d_1 \times R} \\ B \in \mathbb{R}^{d_2 \times R} \\ C \in \mathbb{R}^{d_3 \times R}}} \| T - A \begin{smallmatrix} R \\ \diagdown \\ d_1 \end{smallmatrix} \begin{smallmatrix} R \\ \diagup \\ d_2 \end{smallmatrix} \begin{smallmatrix} R \\ \diagup \\ d_3 \end{smallmatrix} \|_F^2$$

Def: The Frobenius norm of  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ :  $\|T\|_F = \sqrt{\sum_{i,j,k} (T_{ijk})^2}$

2) TUCKER

Given  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  and a target rank  $(R_1, R_2, R_3)$

$$\min_{\substack{G \in \mathbb{R}^{R_1 \times R_2 \times R_3} \\ U_1 \in \mathbb{R}^{d_1 \times R_1} \\ U_2 \in \mathbb{R}^{d_2 \times R_2} \\ U_3 \in \mathbb{R}^{d_3 \times R_3}}} \| T - G \times_{1,2} U_1 \times_2 U_2 \times_3 U_3 \|_F^2$$

3) TENSOR TRAIN

Given  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  and a target rank  $(R_1, R_2)$

$$\min_{\substack{G_1 \in \mathbb{R}^{d_1 \times R_1} \\ G_2 \in \mathbb{R}^{R_1 \times d_2 \times R_2} \\ G_3 \in \mathbb{R}^{R_2 \times d_3}}} \| T - G_1 \begin{smallmatrix} R_1 \\ \diagdown \\ d_1 \end{smallmatrix} \begin{smallmatrix} R_2 \\ \diagup \\ d_2 \end{smallmatrix} G_3 \|_F^2$$

## OPTIMIZATION WITH TENSORS

General problem: Given a loss function  $\mathcal{L}: \mathbb{R}^{d_1 \times \dots \times d_N} \rightarrow \mathbb{R}$ , we want to solve

$$\min_{W \in \mathbb{R}^{d_1 \times \dots \times d_N}} \mathcal{L}(W) \quad \text{subject to } \text{rank}(W) \leq R$$

↳ CP rank, Tucker rank, TT rank, ...

Examples of loss functions:

+ Low rank approximation

$$\min_W \|W - T\|_F^2 \quad \text{s.t. } \text{rank}_{CP}(W) \leq R$$

↳ the loss function is  $\mathcal{L}(W) = \|W - T\|_F^2$

+ Regression We want to learn a linear function  $f: \mathbb{R}^{d_1 \times \dots \times d_N} \rightarrow \mathbb{R}$

from a dataset  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \subseteq \mathbb{R}^{d_1 \times \dots \times d_N \times \mathbb{R}}$ .

$f$  is linear  $\Leftrightarrow f(x) = f_w(x) = \langle w, x \rangle = \sum_{i_1, \dots, i_N} w_{i_1, \dots, i_N} x_{i_1, \dots, i_N}$ .

In this case, a natural loss function is

$$\begin{aligned} \mathcal{L}(W) &= \sum_{i=1}^m (f_w(x_i) - y_i)^2 \\ &= \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \end{aligned}$$

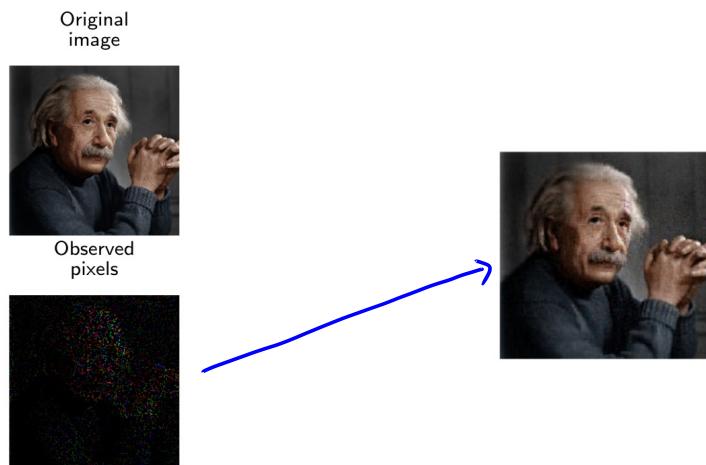
+ Completion: Target tensor  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$

Data: observed entries  $\{T_{i_1, \dots, i_N} | (i_1, \dots, i_N) \in \mathcal{I}\}$

$$\mathcal{I} \subseteq [d_1] \times [d_2] \times \dots \times [d_N]$$

- $\begin{pmatrix} 1 & ? & ? \\ ? & 2 & ? \\ ? & ? & 3 \end{pmatrix} \quad \mathcal{I} = \{(1,1), (2,2), (2,3)\}$

- Image completion



For completion, the loss function is :

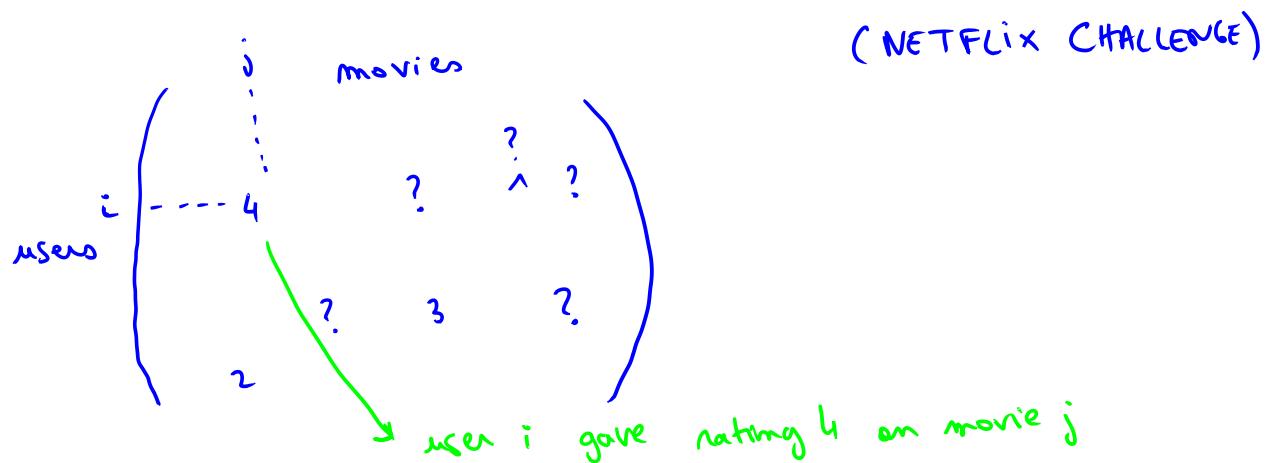
$$\mathcal{L}(W) = \sum_{(i_1, \dots, i_N) \in \mathcal{R}} (T_{i_1, \dots, i_N} - W_{i_1, \dots, i_N})^2$$

*↳ set of observed entries*

Remark: If  $\mathcal{R} = [d_1] \times [d_2] \times \dots \times [d_N]$ , then

$$\sum_{(i_1, \dots, i_N) \in \mathcal{R}} (T_{i_1, \dots, i_N} - W_{i_1, \dots, i_N})^2 = \| T - W \|_F^2$$

- COLLABORATIVE FILTER / RECOMMENDATION SYSTEMS



# ① GENERAL OPTIMIZATION ALGORITHM

**PB1**  $\min_{W \in \mathbb{R}^{d_1 \times \dots \times d_N}} L(W)$  subject to  $\text{rank}_{CP}(W) \leq R$

Notation: Given  $A_i \in \mathbb{R}^{d_i \times R}$  for  $i=1, \dots, N$ , we denote the CP decomposition with factors  $A_1, \dots, A_N$  by  $CP(A_1, \dots, A_N)$

ex:  $CP(A, B, C) = \begin{array}{c} R \\ \overbrace{A \quad B \quad C}^{\text{Tr}} \\ d_1 | \quad d_2 | \quad d_3 | \end{array}$

PB1 is equivalent to:

**PB2**  $\min_{\substack{A_i \in \mathbb{R}^{d_i \times R} \\ i=1, \dots, N}} L(CP(A_1, \dots, A_N))$

## 1) Gradient based algorithm

**ALGORITHM**: Initialize  $A_1, A_2, \dots, A_N$

Repeat

```

    for i=1...N
         $x_i \leftarrow A_i - \gamma \nabla_{A_i} L(CP(A_1, \dots, A_N))$  learning rate
    for i=1...N
         $A_i \leftarrow x_i$ 
until convergence
  
```

## 2) Alternating minimization

We assume  $\min_{A_i} L(CP(A_1, \dots, A_N))$  is easy.

**ALGORITHM**: Initialize  $A_1, A_2, \dots, A_N$

Repeat

```

    for i=1,...,N
         $A_i \leftarrow \underset{A_i}{\text{arg min}} L(CP(A_1, \dots, A_N))$ 
until convergence
  
```

### 3) Alternating minimization for low CP rank approximation

Given  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$ , target rank  $R$ :

$$\boxed{\text{CP-PSM}} \quad \min_{\substack{A_i \in \mathbb{R}^{d_i \times R} \\ i=1, \dots, N}} \| T - \text{CP}(A_1, \dots, A_N) \|_F^2$$

- Matricization:  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$ , for any  $m \in [N]$ ,  $T_{(m)} \in \mathbb{R}^{d_m \times d_1 \cdots d_{m-1} \cdots d_N}$   
 "mode- $m$  matricization"
- Remark:  $\| T \|_F = \| T_{(m)} \|_F$  for all  $m \in [N]$

Def: If  $A \in \mathbb{R}^{m \times m}$  and  $B \in \mathbb{R}^{n \times q}$ , then Kronecker product  
 $A \otimes B \in \mathbb{R}^{mn \times mq}$  is defined by:

$$A \otimes B = \begin{pmatrix} a_{11}B & \cdots & a_{1m}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \cdots & a_{mm}B \end{pmatrix}_{mn \times mq} \quad (\text{a}_{ij} \text{ is the entry } (i,j) \text{ of } A)$$

Remark:  $a \otimes b = \text{vec}(b a^\top) = \text{vec}(\underbrace{b \circ a}_{m \times m})$

Def: The Kathri-Rao product of  $A = \begin{pmatrix} a'_1 & \cdots & a'_R \end{pmatrix} \in \mathbb{R}^{m \times R}$  and  
 $B = \begin{pmatrix} b'_1 & \cdots & b'_R \end{pmatrix} \in \mathbb{R}^{n \times R}$  is defined by

$$A \odot B = \begin{pmatrix} a'_1 \otimes b'_1 & a'_2 \otimes b'_2 & \cdots & a'_R \otimes b'_R \end{pmatrix} \in \mathbb{R}^{mn \times R}$$

Property: If  $W = \text{CP}(A_1, A_2, \dots, A_N)$ , then

$$W_{(m)} = \underbrace{A_m^\top (A_N \odot \cdots \odot A_{m+1} \odot A_{m-1} \odot \cdots \odot A_1)}_{d_N \times d_{m+1} \cdots d_m \times R}^\top$$

## Alternating least squares algorithm (ALS) for CP decomposition

$$\min_{\substack{A, B, C \\ d_1 \times R \\ d_2 \times R \\ d_3 \times R}} \| T - CP(A, B, C) \|_F^2$$

We want solve this problem w.r.t A:

$$\begin{aligned} \underset{A}{\operatorname{argmin}} \| T - CP(A, B, C) \|_F^2 &= \underset{A}{\operatorname{argmin}} \| T_{(1)} - (CP(A, B, C))_{(1)} \|_F^2 \\ &= \underset{A}{\operatorname{argmin}} \| T_{(1)} - A (C \odot B)^T \|_F^2 \\ &= T_{(1)} [ (C \odot B)^T ]^+ \end{aligned}$$

↑  
pseudo-inverse

(side note: if  $A \in \mathbb{R}^{m \times n}$  with  $m \leq n$  and  $\operatorname{rank}(A) = m$ , then

$$AA^+ = I \quad (\text{however: } \Delta A^+ A \text{ may not be } I \Delta)$$

### ALS ALGORITHM:

INPUT: target tensor  $T$ , rank  $R$

OUTPUT: factor matrices  $A_i, i=1, \dots, N$  s.t.  
 $T \approx CP(A_1, \dots, A_N)$

- Initialize  $A_1, \dots, A_N$

- Repeat

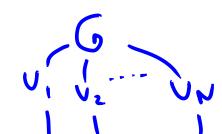
For  $m = 1, \dots, N$   
 $A_m \leftarrow T_{(m)} \left[ (A_N \odot \dots \odot A_{m+1} \odot A_{m-1} \odot \dots \odot A_1)^T \right]^+$   
 until convergence

## II SVD-based algorithms for TUCKER and TT

### 1) TUCKER: Higher order SVD (LATHAWALA et al. 2000)

Given  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$  and  $(R_1, \dots, R_N)$  we want to solve

$$\min_{\substack{G \in \mathbb{R}^{R_1 \times \dots \times R_N} \\ U_i \in \mathbb{R}^{d_i \times R_i}, \\ i=1 \dots N}} \|T - G \times_1 U_1 \times_2 U_2 \times \dots \times_N U_N\|_F^2$$



### HOSVD Algorithm

INPUT:  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$  and  $(R_1, \dots, R_N)$

OUTPUT:  $\begin{array}{l} G \in \mathbb{R}^{R_1 \times \dots \times R_N} \\ U_i \in \mathbb{R}^{d_i \times R_i}, \\ i=1 \dots N \end{array}$  s.t.  $T \approx G \times_1 U_1 \times_2 \dots \times_N U_N$

For  $m = 1 \dots N$

$\left[ \begin{array}{l} U_m \leftarrow R_m \text{ leading left singular vectors of } T_{(m)} \\ d_m \times R_m \end{array} \right]$

$$G \leftarrow T \times_1 U_1^T \times_2 \dots \times_N U_N^T$$

$d_1 \times \dots \times d_N$        $R_1 \times d_1$        $R_N \times d_N$

RETURN  $G, U_1, \dots, U_N$

$\downarrow$  SVD  
 $U D V^T$   
 $\downarrow$   
extract the first  
 $R_m$  columns

★  $\min_X \|T - X\|_F^2 \quad \text{s.t. } \text{rank}_{\text{TUCKER}}(X) \leq (R_1, \dots, R_N)$

Theorem: Let  $X^*$  be the solution of problem ★ and  $X_{\text{HOSVD}}$  the solution returned by HOSVD.

$\hookrightarrow G \times_1 U_1 \times \dots \times_N U_N$

Then

$$\|T - X_{\text{HOSVD}}\|_F \leq \sqrt{N} \|T - X^*\|_F$$

If  $T$  has Tucker rank less than  $(R_1, \dots, R_N)$ , then  $\|T - X^*\|_F = 0$  hence  $\|T - X_{\text{HOSVD}}\|_F = 0$  my HOSVD recover the exact TUCKER decomposition if it exists.

2) TENSOR TRAIN : TT-SVD (OSSEDELETS, 2010)  
 ↳ (a.k.a. MATRIX PRODUCT STATE (MPS) : sequential SVD)

$$\min_{G_1, G_2, \dots, G_N} \| T - \text{TT}(G_1, G_2, \dots, G_N) \|_F^2$$

$$\hookrightarrow \frac{G_1}{d_1} \xrightarrow{R_1} \frac{G_2}{d_2} \xrightarrow{R_2} \frac{G_3}{d_3} \xrightarrow{R_3} \dots \xrightarrow{R_{N-1}} \frac{G_N}{d_N}$$

### TT-SVD Algorithm

INPUT:  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3 \times d_4}$ ,  $(R_1, \dots, R_3)$

OUTPUT:  $G_1, \dots, G_4$  s.t.  $T \approx \text{TT}(G_1, G_2, G_3, G_4)$

- $\frac{d_1}{d_1} \frac{T}{\overset{d_2}{\underset{d_4}{\diagup}} \underset{d_3}{\diagdown}} \approx \frac{d_2}{d_1} G_1 \xrightarrow{R_1} A \underset{d_4}{\overset{d_3}{\equiv}} \frac{d_3}{d_4}$  (Low rank matrix approx.  $\rightarrow$  SVD)
- $\frac{d_2}{R_1} A \underset{d_4}{\overset{d_3}{\equiv}} \approx \frac{d_2}{R_1} G_2 \xrightarrow{R_2} B \underset{d_4}{\overset{d_3}{\equiv}}$
- $\frac{d_3}{R_2} B \underset{d_4}{\overset{d_4}{\equiv}} \approx \frac{d_3}{R_2} G_3 \xrightarrow{R_3} G_4 \frac{d_4}{d_4}$
- RETURN  $G_1, G_2, G_3, G_4$

Theorem: If  $X^*$  is the best approximation of  $T$  of TT rank  $(R_1, \dots, R_{N-1})$  then  $\|T - X_{\text{TT-SVD}}\|_F \leq \sqrt{N-1} \|T - X^*\|_F$ .

↳ solution returned by TT-SVD

$$X_{\text{TT-SVD}} = \frac{G_1}{1} - \frac{G_2}{1} - \dots - \frac{G_N}{1}$$

- If  $\|T - X^*\|_F = 0$  then TT-SVD returns an exact TT decomposition of  $T$ .

### (III) SUPERVISED LEARNING With TENSOR NETWORKS

Supervised Learning With Quantum-Inspired Tensor Networks

(NeurIPS, 2016)

E. Miles Stoudenmire<sup>1,2</sup> and David J. Schwab<sup>3</sup>

We want to learn a linear function  $f : \mathbb{R}^{d^m} \rightarrow \mathbb{R}^n$   
 input space is very high dimension  
 $\downarrow$   
 $f: x \mapsto Wx$   
 $W \in \mathbb{R}^{n \times d^m}$

Feature map: Image  $\rightarrow \mathbb{R}^{d \times d \times \dots \times d}$

$x \in \mathbb{R}^{h \times w} \mapsto \underbrace{\phi(x_{11}) \circ \phi(x_{12}) \circ \phi(x_{13}) \circ \dots \circ \phi(x_{hw})}_{\text{Tensor of order } hw} \rightarrow \mathbb{R}^{2^{hw}}$

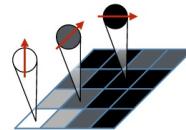
Two examples for  $\phi$  are

$$\begin{aligned} \cdot \phi : \mathbb{R} &\rightarrow \mathbb{R}^2 \\ \alpha &\mapsto \begin{pmatrix} \alpha \\ 1 \end{pmatrix} \end{aligned}$$

$$\hookrightarrow (d_1, d_2, d_3) \mapsto \underbrace{\phi(d_1) \circ \phi(d_2) \circ \phi(d_3)}_{\text{Entries of this tensor}} = \begin{pmatrix} d_1 \\ 1 \end{pmatrix} \circ \begin{pmatrix} d_2 \\ 1 \end{pmatrix} \circ \begin{pmatrix} d_3 \\ 1 \end{pmatrix}$$

Entries of this tensor:  $1, d_1, d_2, d_3, d_1 d_2, d_1 d_3, d_2 d_3$   
 and  $d_1 d_2 d_3$

$$\begin{aligned} \cdot \phi : \mathbb{R} &\rightarrow \mathbb{R}^2 \\ \alpha &\mapsto \begin{pmatrix} \cos\left(\frac{\pi}{2}\alpha\right) \\ \sin\left(\frac{\pi}{2}\alpha\right) \end{pmatrix} \end{aligned}$$



$$f: \mathbb{R}^{d^m} \approx \mathbb{R}^{d \times d \times \dots \times d} \xrightarrow{m} \mathbb{R}^n$$

$\xrightarrow{\quad}$

$$\begin{matrix} X \\ d/d \dots /d \end{matrix} \xrightarrow{\quad} \begin{matrix} X \\ d(d \dots) d \\ \diagdown \quad \diagup \\ W \\ | n \end{matrix}$$

$\hookrightarrow$  The model parameter is  $W \in \mathbb{R}^{d \times \dots \times d \times n}$   
 $\hookrightarrow d^m n$  parameters!

We cannot even store  $W$  in memory!

Sol: Parameterize  $W$  with a low rank TT decomposition

$$\begin{matrix} | n \\ W \\ d/d \dots /d \end{matrix} = \begin{matrix} O & O & O & O & \dots & O & | n \\ d & d & d & d & \dots & d \end{matrix}$$

$$f(X) = f(\phi(x_{11}) \circ \phi(x_{12}) \circ \dots \circ \phi(x_{1w}))$$

$$= \begin{matrix} \phi(x_{11}) & \phi(x_{12}) & \dots & \phi(x_{1w}) \\ d | & d | & \dots & | d \\ \underbrace{\quad}_{W} \end{matrix}$$

$$= \begin{matrix} \phi(x_{11}) & \phi(x_{12}) & \dots & \phi(x_{1w}) \\ d | & d | & \dots & | d \\ G_1 - G_2 - \dots - G_{nw} - G_{nw+1} & | n \end{matrix} \quad \begin{matrix} \text{We can compute} \\ \text{this very} \\ \text{efficiently} \end{matrix}$$

Learning problem: Given data  $\{(x_1, y_1), \dots, (x_m, y_m)\} \subseteq \mathbb{R}^{d^m} \times \mathbb{R}^n$

we want to minimize  $\hat{L} = \sum_{i=1}^m l(f(x_i), y_i)$  for some loss function

$$l: \mathbb{R}^n \rightarrow \mathbb{R}.$$

$\hookrightarrow$  parameters:  $G_1, G_2, \dots, G_{nw+1}$  cores of a TT decomposition of  $W$ .

## Training Algorithm (DMRG)

. Initialize  $G_1, G_2, \dots, G_{m+1}$   
 . Repeat  
 | for each consecutive pair of cores  $G_i, G_{i+1}$   
 | (merge two cores)  
 | (Gradient descent step)  
 | (split in two cores)  

$$\overbrace{B}^{R_{i-1} R_{i+1}} = \overbrace{G_i}^{d \times d} \overbrace{G_{i+1}}^{d \times d}$$
  

$$B^{\text{new}} = B - \gamma \nabla_B L$$
  

$$\overbrace{B}^{\text{new}} \approx \overbrace{G_i}^{d \times d} \overbrace{R_i}^{\text{new}} \overbrace{G_{i+1}}^{\text{new}} \overbrace{R_{i+1}}^{\text{new}}$$
 (truncated SVD)  

$$G_i, G_{i+1} \leftarrow G_i^{\text{new}}, G_{i+1}^{\text{new}}$$
  
 until convergence  
 \*  $R_i^{\text{new}}$  can be chosen adaptively from the singular values of  $\overbrace{B}^{\text{new}}$ .

Follow up notes:

From probabilistic graphical models to generalized tensor networks  
for supervised learning

Ivan Glasser,<sup>1,2</sup> Nicola Pancotti,<sup>1,2</sup> and J. Ignacio Cirac<sup>1,2</sup>

### Tensor Networks for Probabilistic Sequence Modeling

Jacob Miller  
Mila and DIRO  
Université de Montréal  
Montréal QC, Canada  
jmjacobmiller@gmail.com

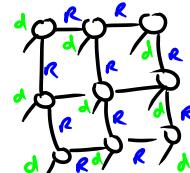
Guillaume Rabusseau  
Mila and DIRO  
Université de Montréal  
Montréal QC, Canada  
grabus@iro.umontreal.ca

John Terilla  
CUNY and Tunnel  
City University of New York  
New York NY, USA  
jterilla@gc.cuny.edu

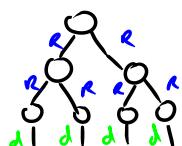
~~TV~~ BEYOND TUCKER, TT, ...

TT : 

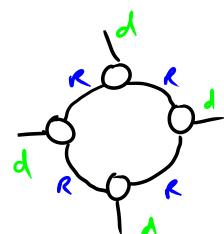
PEPS :  
(2dim TN)



HIERARCHICAL TUCKER :  
(Tree TN)



TENSOR RING :



# TENSOR METHODS FOR ML

- COMPRESS NEURAL NETWORKS
- LEARNING LATENT VARIABLE MODELS

## ① EFFICIENT COMPUTATIONS WITH TENSOR NETWORKS

- Inner product:  $u, v \in \mathbb{R}^{d^N}$

What is the computational cost for  $\langle u, v \rangle = \sum_{i=1}^{d^N} u_i v_i$ ?  
 $\hookrightarrow \mathcal{O}(d^N)$

Suppose  $u$  and  $v$  are in the TT format:

$$u, v \in \mathbb{R}^{d^N} \approx \mathbb{R}^{\underbrace{d \times d \times \dots \times d}_{N \text{ times}}}$$

$$\begin{aligned} u &= \underbrace{U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \\ v &= \underbrace{V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \end{aligned} \quad \left. \right\} \text{TT representations of } u \text{ and } v$$

$$\langle u, v \rangle = u - v = \underbrace{U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \quad \underbrace{V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow}$$

Recall: Computing the matrix product  $A B$  has cost  $\mathcal{O}(m n p)$

$$(d^2 \times R) \quad (R \times d^2) \quad m \times n \times p$$

For ex, Computing  $\underbrace{A \xrightarrow{R}}_{d \times 1}, \underbrace{B \xrightarrow{1 \times d}}_{1 \times d} \rightsquigarrow \underbrace{A \xrightarrow{R} B \xrightarrow{1 \times d}}_{d \times 1}$  can be done

with a matrix product between matrices of shapes  $d^2 \times R$  and  $R \times d^2$ . Hence the complexity is  $\mathcal{O}(R d^4)$ .

We want to compute

$$U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N$$

$$V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N$$

### COMPLEXITY

① Compute  $U_1 \xrightarrow{R}$

$$\begin{array}{c} d \\ | \\ U_1 \xrightarrow{R} \end{array}$$

$\mathcal{O}(dR^2)$

$$\begin{array}{c} R \times R \\ | \\ U_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ U_2 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} d \\ | \\ V_1 \xrightarrow{R} \end{array}$$

$\mathcal{O}(dR^3)$

③  $U_1 \xrightarrow{R} U_2 \xrightarrow{R}$

$$\begin{array}{c} d \\ | \\ U_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} d \\ | \\ U_2 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} V_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} V_2 \xrightarrow{R} \\ | \\ R \times dR \end{array}$$

$\mathcal{O}(dR^3)$

$$\begin{array}{c} R \times R \\ | \\ U_1 \xrightarrow{R} U_2 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ U_3 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} d \\ | \\ V_1 \xrightarrow{R} V_2 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ R \times dR \end{array}$$

$\mathcal{O}(dR^3)$

⑤

---

:

:

$$U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N$$

$$V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N$$

this is linear in  $N$

TOTAL COMPLEXITY

$\mathcal{O}(NdR^3)$

↳ in contrast with  $\mathcal{O}(d^N)$  !  
 exp. in  $N$

- Matrix vector product  
 $A \in \mathbb{R}^{d^N \times n^N}$ ,  $x \in \mathbb{R}^{n^N}$  : we want to compute  $b = Ax \in \mathbb{R}^{d^N}$ .

In fact, we want to compute  $b$  in the TT format.

+ TT representation of matrices :

$$A \in \mathbb{R}^{d^N \times n^N} \approx \mathbb{R}^{\underbrace{d \times d \times \dots \times d}_{N \text{ times}} \times \underbrace{n \times n \times \dots \times n}_{N \text{ times}}}$$

① First solution :

$$A = A_1 \frac{R}{d} A_2 \frac{R}{d} \dots \frac{R}{d} A_N \frac{R}{n} A_{N+1} \frac{R}{n} A_{N+2} \frac{R}{n} \dots \frac{R}{n} A_{2N}$$

↳ PBM :  $A$  will be of rank at most  $R$ .

$$A = \begin{matrix} A_1 \frac{R}{d} A_2 \frac{R}{d} \dots \frac{R}{d} A_N \end{matrix} \frac{R}{n} \begin{matrix} A_{N+1} \frac{R}{n} A_{N+2} \frac{R}{n} \dots \frac{R}{n} A_{2N} \end{matrix}$$

$P$                                      $Q$

$d^N \times R$                              $R \times n^N$

$$\hookrightarrow A = \begin{matrix} d^N \times n^N \\ / \quad \backslash \end{matrix} P Q \begin{matrix} R \times n^N \end{matrix} \rightarrow \text{Rank } R \text{ factorization of } A, \text{ hence } \text{rank}(A) \leq R.$$

② TT matrix representation :

$$A = \begin{matrix} d \\ -A_1 \frac{n}{R} \\ d \\ -A_2 \frac{n}{R} \\ d \\ -A_3 \frac{n}{R} \\ \vdots \\ d \\ -A_N \frac{n}{R} \end{matrix}$$

↳ Remark • If  $A = \begin{matrix} d^2 \times n^2 \\ -A_1 \frac{n}{R} \\ -A_2 \frac{n}{R} \end{matrix}$  and  $R = 1$ ,

$$\text{then } A = \begin{matrix} d^2 \times n^2 \\ -A_1 \frac{n}{R} \\ -A_2 \frac{n}{R} \end{matrix} = A_1 \otimes A_2 \quad \text{KRONECKER PRODUCT}$$

• In a TN, an edge of dimension 1 is equivalent to having no edge :

$$\left( \begin{matrix} -B & C \\ m \times 1 & 1 \times m \end{matrix} \right)_{ij} = \sum_{k=1}^1 B_{ik} C_{kj} = B_{i1} C_{j1} = b \circ c$$

vectors       $b$        $c^T$

Back to matrix vector product:

$$A = \begin{array}{c} d \\ d^4 \times h^4 \end{array} = \begin{array}{c} d \\ d \\ d \\ d \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array}$$

$$x = \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$$

We want a TT representation of  $b = Ax$

	$d \times R \times R$	$d \times R^2$	COMPLEXITY
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_1 \\ B_2 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^2 d n)$
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_2 \\ B_3 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^4 d n)$
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} B_1 \\ B_2 \\ B_3 \\ B_4 \end{array} \begin{array}{c} R^2 \\ R^2 \\ R^2 \\ R^2 \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_3 \\ B_4 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^4 d n)$
			$O(R^2 d n)$

$\Rightarrow$  If  $A \in \mathbb{R}^{d^N \times h^N}$  given as a TT matrix with rank  $R$

and  $x \in \mathbb{R}^{h^N}$  given as a TT vector with rank  $R$

then a rank  $R^2$  TT representation of  $b = Ax$  can be computed in time  $O(N R^4 d n)$ .

↳ in contrast with  $O(d^N h^N)$

- $u, v \in \mathbb{R}^{d^N}$ , given as rank  $R$  TT-vector, a rank  $2R$  TT representation of  $u + v$  can be computed in  $\mathcal{O}(NR^3d)$ .
- $u, v \in \mathbb{R}^{d^N}$ , given as rank  $R$  TT-vector, a rank  $R^2$  TT representation of  $u \otimes v$  can be computed in  $\mathcal{O}(NR^3d)$ .  
↳ component-wise product.

- TT Rounding:

Given a rank  $R$  representation of  $u \in \mathbb{R}^{d^N}$ , the TT rounding operation returns a rank  $\hat{R}$  TT representation of  $u$  which is "close" to the original representation (in time  $\mathcal{O}(NR^3d)$ ).

FOR MORE DETAILS:

SIAM J. SCI. COMPUT.  
Vol. 33, No. 5, pp. 2295–2317

© 2011 Society for Industrial and Applied Mathematics

## TENSOR-TRAIN DECOMPOSITION\*

I. V. OSELEDETS†

**Abstract.** A simple nonrecursive form of the tensor decomposition in  $d$  dimensions is presented. It does not inherently suffer from the curse of dimensionality; it has asymptotically the same number of parameters as the canonical decomposition, but it is stable and its computation is based on low-rank approximation of auxiliary *unfolding matrices*. The new form gives a clear and convenient way to implement all basic operations efficiently. A fast rounding procedure is presented, as well as basic linear algebra operations. Examples showing the benefits of the decomposition are given, and the efficiency is demonstrated by the computation of the smallest eigenvalue of a 19-dimensional operator.

## II COMPRESSING NEURAL NETWORKS

### Tensorizing Neural Networks

(NeurIPS, 2015)

Alexander Novikov<sup>1,4</sup> Dmitry Podoprikin<sup>1</sup> Anton Osokin<sup>2</sup> Dmitry Vetrov<sup>1,3</sup>

<sup>1</sup>Skolkovo Institute of Science and Technology, Moscow, Russia

<sup>2</sup>INRIA, SIERRA project-team, Paris, France

<sup>3</sup>National Research University Higher School of Economics, Moscow, Russia

<sup>4</sup>Institute of Numerical Mathematics of the Russian Academy of Sciences, Moscow, Russia

novikov@bayesgroup.ru podoprikin.dmitry@gmail.com

anton.osokin@inria.fr vetrov@yandex.ru

Fully connected layer:

size:  $M \times N$

$$y = \sigma(Wx + b)$$

activation function

For sake of simplicity, we ignore the bias term:  $y = \sigma(Wx)$

IDEA: Represent  $W$  as a TT matrix.

$$M = m_1 \times m_2 \times \dots \times m_d$$

( $d$  is the order of the tensor)

$$N = n_1 \times n_2 \times \dots \times n_d$$

$$W = \begin{matrix} & m_1 & m_1 \\ & G_1 & | \\ m_2 & G_2 & m_2 \\ & | & | \\ & \vdots & \vdots \\ & G_d & m_d \end{matrix}$$

TT layer

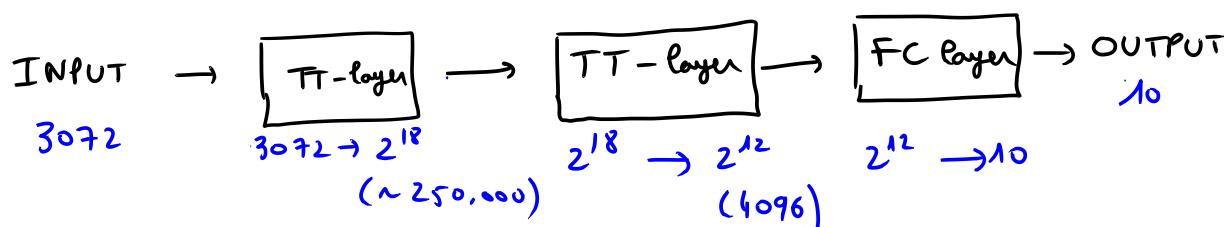
↳ parameters are  
 $G_1, G_2, \dots, G_d$

COMPLEXITY:

Operation	Time	Memory
FC forward pass	$O(MN)$	$O(MN)$
TT forward pass	$O(dr^2 m \max\{M, N\})$	$O(r \max\{M, N\})$
FC backward pass	$O(MN)$	$O(MN)$
TT backward pass	$O(d^2 r^4 m \max\{M, N\})$	$O(r^3 \max\{M, N\})$

Table 1: Comparison of the asymptotic complexity and memory usage of an  $M \times N$  TT-layer and an  $M \times N$  fully-connected layer (FC). The input and output tensor shapes are  $m_1 \times \dots \times m_d$  and  $n_1 \times \dots \times n_d$  respectively ( $m = \max_{k=1 \dots d} m_k$ ) and  $r$  is the maximal TT-rank.

EXPERIMENT ON CIFAR 10:



COMPUTING GRADIENTS OF TENSOR NETWORKS

↳ See TN-gradients.hdf



## Tensor Decompositions for Learning Latent Variable Models

**Animashree Anandkumar**  
*Electrical Engineering and Computer Science  
 University of California, Irvine  
 2200 Engineering Hall  
 Irvine, CA 92697*

A.ANANDKUMAR@UCI.EDU

(JMLR, 2014)

**Rong Ge**  
*Microsoft Research  
 One Memorial Drive  
 Cambridge, MA 02142*

RONGGE@MICROSOFT.COM

**Daniel Hsu**  
*Department of Computer Science  
 Columbia University  
 1214 Amsterdam Avenue, #0401  
 New York, NY 10027*

DJHSU@CS.COLUMBIA.EDU

**Sham M. Kakade**  
*Microsoft Research  
 One Memorial Drive  
 Cambridge, MA 02142*

SKAKADE@MICROSOFT.COM

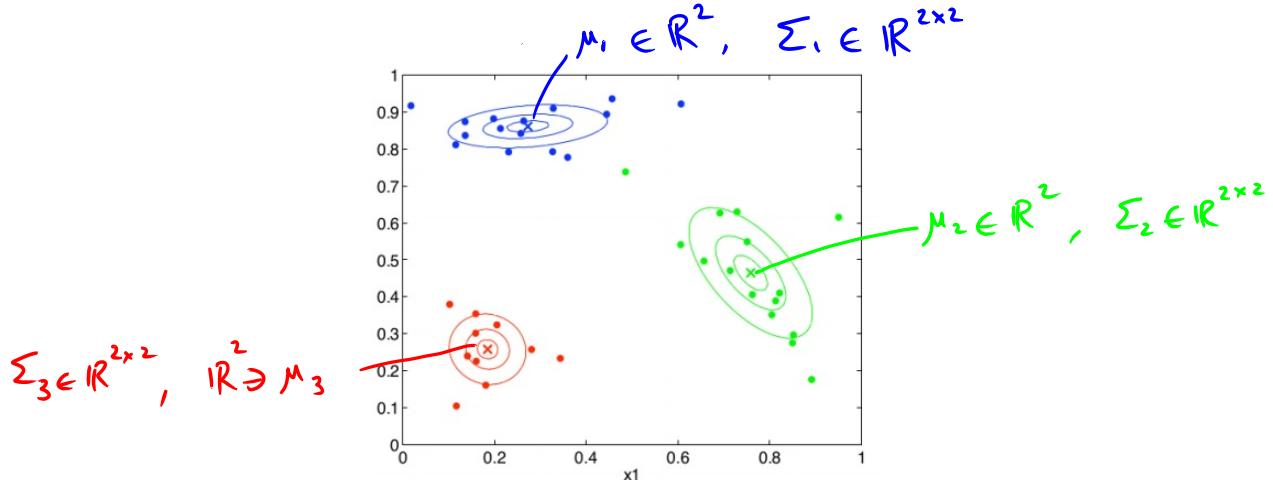
**Matus Telgarsky**  
*Department of Statistics  
 Rutgers University  
 110 Frelinghuysen Road  
 Piscataway, NJ 08854*

MTELGARS@CS.UCSD.EDU

Editor: Benjamin Recht

## 1) Examples of LVM

## • Gaussian mixture model (GMM)

GMM with  $k$  components in  $\mathbb{R}^d$ :• parameters :  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$  (centers)n.s.d. matrices  $\Sigma_1, \dots, \Sigma_k \in \mathbb{R}^{d \times d}$  (covariance matrices) $0 \leq \pi_i \leq 1, \sum_{i=1}^k \pi_i = 1 \leftarrow \pi_1, \dots, \pi_k \in \mathbb{R}_+$  (mixing probabilities)

• data generating distribution

Draw  $x \sim \text{GMM}$

(i) Draw a Gaussian  $h$  with probabilities

$$\Pr(h=i) = \pi_i, i=1\dots,k$$

(ii) Draw  $x \sim \mathcal{N}(\mu_h, \Sigma_h)$

• Single topic model

(latent variable is the topic of a document)

Topic: discrete random variable taking values  $1, 2, \dots, k$

Vocabulary of size  $d$  ( $d$  different words).

→ We want to model the distribution over the vocabulary given a topic.

• Parameters:

$$P(\text{topic} = i) = \pi_i \quad \pi_1, \pi_2, \dots, \pi_k \in \mathbb{R} \quad (\text{topic probabilities})$$

$$0 \leq \pi_i \leq 1, \sum_{i=1}^k \pi_i = 1$$

$$P(\text{word} = j | \text{topic} = i) = (\mu_i)_j \quad \mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d \quad (\text{voc. distrib. given a topic})$$

$$\text{for } i=1, \dots, k : \quad 0 \leq (\mu_i)_j \leq 1, \sum_{j=1}^d (\mu_i)_j = 1$$

• Data generating distribution

Draw a document of length  $\ell$  from the single topic model.

Draw a topic $h$ with probabilities $P(h = i) = \pi_i$	• Draw independently $\ell$ words with probabilities $P(\text{word } j   \text{topic} = h) = (\mu_h)_j$
---	--

## 2) Moments and Latent variable models

Def: • If  $x$  is a random variable taking its values in  $\mathbb{R}$ , its  $n^{\text{th}}$  order moment is  $\mathbb{E}[x^n]$ .

• If  $x$  is a random variable taking its values in  $\mathbb{R}^d$ , its  $n^{\text{th}}$  order moment is  $\mathbb{E}[x^{0:n}] = \mathbb{E}\left[\underbrace{x_0 x_0 \dots x_0}_{n \text{ times}}\right] \in \mathbb{R}^{\frac{d \times d \times \dots \times d}{n \text{ times}}}$ .

↳ . 3<sup>rd</sup> moment :  $\mathbb{E}[x_0 x_0 x_0]_{i,j,k} = \mathbb{E}[x_i x_j x_k]$

. 1<sup>st</sup> moment :  $\mathbb{E}[x] \rightarrow \text{mean}$

. 2<sup>nd</sup> moment :  $\mathbb{E}[x_0 x] = \mathbb{E}[xx^T]$

↳ if  $\mathbb{E}[x] = 0$ , this is the covariance matrix.

↳  $\mathbb{V}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$

$$\text{Cov}[x] = \mathbb{E}[xx^T] - \mathbb{E}[x]\mathbb{E}[x]^T$$

→ the method of moments give the same estimation as the maximum likelihood for a Gaussian distribution.

$$\text{If } x \sim \mathcal{N}(\mu, \sigma^2) \quad \mathbb{E}[x] = \mu \quad \mathbb{E}[x^2] = \sigma^2 + \mu^2 \quad \left| \begin{array}{l} \mu = \mathbb{E}[x] \approx \overbrace{\frac{1}{m} \sum_{i=1}^m x_i}^{\hat{\mu}} \\ \sigma^2 = \mu^2 - \mathbb{E}[x^2] \approx \hat{\mu}^2 - \overbrace{\frac{1}{m} \sum_{i=1}^m x_i^2}^{\hat{\sigma}^2} \end{array} \right.$$

### Simple Topic Model

$$\text{parameters: } \Theta = \{\mu_1, \mu_2, \dots, \mu_K, \mu_{1,1}, \dots, \mu_{K,K}\}$$

$\downarrow \mathbb{R}$        $\downarrow \mathbb{R}^d$

$$\begin{aligned} \mathbb{P}(\text{word } i, j \mid \text{topic } h) &= \mathbb{P}(\text{word } i \mid \text{topic } h) \mathbb{P}(\text{word } j \mid \text{topic } h) \\ &= (\mu_h)_i (\mu_h)_j \\ &= (\mu_h \circ \mu_h)_{i,j} \end{aligned}$$

$$\begin{aligned} \mathbb{P}(\text{word } i, j) &= \sum_{h=1}^K \mathbb{P}(\text{topic } h) \mathbb{P}(\text{word } i, j \mid \text{topic } h) \\ &= \sum_{h=1}^K \mu_h (\mu_h \circ \mu_h)_{i,j} \end{aligned}$$

$$\mathbb{P}(\text{word } i_1, i_2, i_3) = \sum_{h=1}^K \mu_h (\mu_h \circ \mu_h \circ \mu_h)_{i_1, i_2, i_3}$$

Note that if  $x_1, x_2$  and  $x_3$  are the one-hot encodings of the first 3 words in a document, then

$$\mathbb{E}[x_1 \circ x_2]_{i_1, i_2} = \mathbb{P}(\text{word } i_1, i_2)$$

$$\mathbb{E}[x_1 \circ x_2 \circ x_3]_{i_1, i_2, i_3} = \mathbb{P}(\text{word } i_1, i_2, i_3)$$

Theorem 3.1 (Anandkumar et al., 2012c) If

$$\begin{aligned} M_2 &:= \mathbb{E}[x_1 \otimes x_2] \\ M_3 &:= \mathbb{E}[x_1 \otimes x_2 \otimes x_3], \end{aligned}$$

then

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned} \quad (\text{w}_i = \mu_i)$$

GMM

parameters :  $\Theta = \left\{ \underbrace{\mu_1, \dots, \mu_k}_{\substack{\text{centers} \\ \mathbb{R}^d}}, \underbrace{\Sigma_1, \dots, \Sigma_k}_{\substack{\text{Covariance} \\ \mathbb{R}^{d \times d}}}, \underbrace{w_1, \dots, w_k}_{\text{mixing probabilities}} \right\}$

↳ Here we assume each  $\Sigma_i = \sigma_i^2 I$  for each  $i=1, \dots, k$   
 ↳  $\sigma_i \in \mathbb{R}$

Theorem 3.2 (Hsu and Kakade, 2013) Assume  $d \geq k$ . The variance  $\sigma^2$  is the smallest eigenvalue of the covariance matrix  $\mathbb{E}[x \otimes x] - \mathbb{E}[x] \otimes \mathbb{E}[x]$ . Furthermore, if

$$\begin{aligned} M_2 &:= \mathbb{E}[x \otimes x] - \sigma^2 I \\ M_3 &:= \mathbb{E}[x \otimes x \otimes x] - \sigma^2 \sum_{i=1}^d (\mathbb{E}[x] \otimes e_i \otimes e_i + e_i \otimes \mathbb{E}[x] \otimes e_i + e_i \otimes e_i \otimes \mathbb{E}[x]), \end{aligned}$$

then

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned}$$

$$\sigma = \sigma_1 = \sigma_2 = \dots = \sigma_k$$

( $e_i$  : i-th vector of the canonical basis)

Summary : In both cases, there exists  $M_2 \in \mathbb{R}^{d \times d}$  and a tensor  $M_3 \in \mathbb{R}^{d \times d \times d}$   
 (which can be estimated from data) satisfying

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned} \quad \rightarrow \text{CP decomposition of } M_3 \text{ (under mild conditions, this decomposition unique)}$$

## Decomposition of $M_2$ and $M_3$ :

$$M_2 = \sum_{i=1}^k w_i \mu_i \otimes \mu_i$$

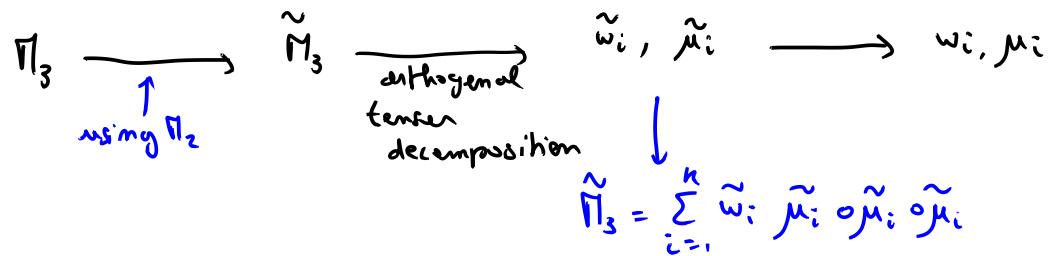
$$M_3 = \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i.$$

Observation: . The problem of orthogonal tensor decomposition,

$$T = \sum_{i=1}^k w_i \mu_i \circ \mu_i \circ \mu_i \quad \text{where } \langle \mu_i, \mu_j \rangle = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{o.w.} \end{cases}$$

is easy to solve.

- We can use  $M_2$  to "orthogonalize"  $M_3$ .



# COMPUTING GRADIENTS OF TENSOR NETWORKS

## I Jacobian and backpropagation

- $f: \mathbb{R}^m \rightarrow \mathbb{R}$       Gradient       $\nabla_{\theta} f = \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_m} \end{pmatrix}$  for each  $\theta \in \mathbb{R}^m$

- $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$       Jacobian       $\frac{\partial f}{\partial \theta} = \left( \frac{\partial f(\theta)_i}{\partial \theta_j} \right)_{i,j} \in \mathbb{R}^{n \times m}$  for each  $\theta \in \mathbb{R}^m$

- Backpropagation / Automatic Differentiation (reverse mode).

Computational Graph       $\theta \rightarrow T^1 \rightarrow T^2 \rightarrow T^3 \rightarrow \mathcal{L}$  (where  $X \rightarrow Y$  means that  $Y$  is a function of  $X$ )  
 $\mathbb{R}^m \rightarrow \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_3} \rightarrow \mathbb{R}$

We want to compute  $\frac{\partial \mathcal{L}}{\partial \theta}$ .

Chain rule:  $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial T_3} \frac{\partial T_3}{\partial T_2} \frac{\partial T_2}{\partial T_1} \frac{\partial T_1}{\partial \theta}$

$\overline{\theta}$        $1 \times m$        $1 \times d_3$        $d_3 \times d_2$        $d_2 \times d_1$        $d_1 \times m$

We define the adjoint:  $\overline{T} = \frac{\partial \mathcal{L}}{\partial T}$

$$\begin{aligned} \overline{T}_3 &= \frac{\partial \mathcal{L}}{\partial T_3} & \overline{T}_2 &= \frac{\partial \mathcal{L}}{\partial T_2} = \frac{\partial \mathcal{L}}{\partial \overline{T}_3} \frac{\partial \overline{T}_3}{\partial T_2} = \overline{T}_3 \frac{\partial T_3}{\partial T_2} & \overline{T}_1 &= \frac{\partial \mathcal{L}}{\partial T_1} = \frac{\partial \mathcal{L}}{\partial T_2} \frac{\partial T_2}{\partial T_1} = \overline{T}_2 \frac{\partial T_2}{\partial T_1} \\ \overline{\theta} &= \dots = \overline{T}_1 \frac{\partial T_1}{\partial \theta} \end{aligned} \quad \boxed{\overline{T}_j = \overline{T}_{j+1} \frac{\partial T_{j+1}}{\partial \overline{T}_j}}$$

More complex graphs:  $\theta \rightarrow T_1 \xrightarrow{T_2} T_4 \rightarrow L$

$$\overline{T_i} = \sum_{j: \text{child of } i} \overline{T_j} \frac{\partial T_j}{\partial T_i}$$

$$\begin{aligned} \text{ex: } \overline{T_1} &= \frac{\partial L}{\partial T_1} = \frac{\partial L}{\partial T_2} \frac{\partial T_2}{\partial T_1} + \frac{\partial L}{\partial T_3} \frac{\partial T_3}{\partial T_1} \\ &= \overline{T_2} \frac{\partial T_2}{\partial T_1} + \overline{T_3} \frac{\partial T_3}{\partial T_1} \end{aligned}$$

$f: \mathbb{R}^{m_1 \times m_2 \times \dots \times m_N} \rightarrow \mathbb{R}^{k_1 \times k_2 \times \dots \times k_p} \quad (\cong \mathbb{R}^{m_1 m_2 \dots m_N} \rightarrow \mathbb{R}^{k_1 k_2 \dots k_p})$

Jacobian Tensor  $\frac{\partial f(T)}{\partial T} = \left( \frac{\partial f(T)_{i_1 \dots i_p}}{\partial T_{j_1 \dots j_N}} \right)_{i_1 \dots i_p, j_1 \dots j_N} \in \mathbb{R}^{k_1 \times \dots \times k_p \times m_1 \times \dots \times m_N}$

## II Jacobian of Tensor Networks

Let  $W$  be a tensor given as a tensor network where  $G$  is a core tensor appearing only once. Then  $\frac{\partial W}{\partial G}$  is simply obtained by deleting  $G$  in the tensor network.

Ex:  $W = \begin{array}{c} d_1 \quad R \\ \diagup \quad \diagdown \\ A \quad B \quad D \\ \diagdown \quad \diagup \\ d_2 \quad C \quad R_2 \\ \diagup \quad \diagdown \\ d_3 \quad R_3 \\ \diagup \quad \diagdown \\ d_4 \end{array}$

$$\frac{\partial W}{\partial B} = \begin{array}{c} d_1 \quad R \\ \diagup \quad \diagdown \\ A \quad D \\ \diagdown \quad \diagup \\ d_2 \quad R_2 \\ \diagup \quad \diagdown \\ d_3 \\ \diagup \quad \diagdown \\ d_4 \end{array}$$

$d_1 \times d_2 \times d_3 \times d_4 \times R_1 \times R_2 \times R_3$

$$\frac{\partial W}{\partial D} = \begin{array}{c} d_1 \quad R_3 \\ \diagup \quad \diagdown \\ A \quad B \\ \diagup \quad \diagdown \\ C \\ \diagup \quad \diagdown \\ d_2 \\ \diagup \quad \diagdown \\ d_3 \\ \diagup \quad \diagdown \\ d_4 \end{array}$$

$d_1 \times d_2 \times d_3 \times d_4 \times R_3 \times d_3 \times d_4$

• Rederiving classical matrix/vector gradients:

$$\frac{\partial \langle u, v \rangle}{\partial u} = \frac{\partial (u \cdot v)}{\partial u} = -v = v$$

$$\frac{\partial Ax}{\partial x} = \frac{\partial (-A \cdot x)}{\partial x} = -A = A$$

$$\frac{\partial x^T Ax}{\partial A} = \frac{\partial (x \cdot A \cdot x)}{\partial A} = x \cdot -x = xx^T$$

$$\frac{\partial T_n(A)}{\partial A} = \frac{\partial (\overset{\circ}{A})}{\partial A} = I = I$$

$$\frac{\partial Ax}{\partial A} = \frac{\partial (-A \cdot x)}{\partial A} = -x \cong I \circ x$$

If the core tensor  $G$  appears  $k$  times in the tensor network of  $W$ , then  $\frac{\partial W}{\partial G}$  is given the sum of  $k$  copies of  $W$  where a different occurrence of  $G$  is deleted in each copy.

$$\underline{Ex:} \quad \frac{\partial x^T A x}{\partial x} = \frac{\partial (x - A^T x)}{\partial x} = -A_x + x - A_x = (A + A^T)x$$

### Applications

TENSORIZING  
NEURAL  
NETWORKS

$$y = \begin{matrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{matrix} \times \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{matrix} \times \begin{matrix} R \\ R \\ R \\ R \end{matrix} \times \begin{matrix} W \\ W \\ W \\ W \end{matrix} \times \begin{matrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{matrix}$$

$$W =$$

$$\begin{matrix} h_1 & d_1 \\ -G_1- \\ h_2 & d_2 \\ -G_2- \\ h_3 & d_3 \\ -G_3- \\ h_4 & d_4 \\ -G_4- \end{matrix}$$

Some loss depending on  $W$

↳ We want to compute  $\frac{\partial L}{\partial G_2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial G_2}$

We need to compute  $\frac{\partial y}{\partial G_2} \in \mathbb{R}^{h_1 \times h_2 \times h_3 \times h_4 \times R \times h_2 \times d_2 \times R}$

$$\frac{\partial y}{\partial G_2} = \left( \begin{matrix} -G_1 \\ -G_2 \\ -G_3 \\ -G_4 \end{matrix} \right) \times \left( \begin{matrix} R \\ R \\ R \\ R \end{matrix} \right)$$



We have:

$$\begin{matrix} h_1 & d_1 \\ -G_1- \\ h_2 & d_2 \\ -G_2- \\ h_3 & d_3 \\ -G_3- \\ h_4 & d_4 \\ -G_4- \end{matrix} \times \begin{matrix} R \\ R \\ R \\ R \end{matrix}$$

## IV Generalized CP decomposition

GENERALIZED CANONICAL POLYADIC TENSOR DECOMPOSITION\*

DAVID HONG<sup>†</sup>, TAMARA G. KOLDA<sup>‡</sup>, AND JED A. DUERSCH<sup>§</sup>

↳ Likelihood maximization & generalized linear model.

Observations/Data :  $X \in \mathbb{R}^{m_1 \times \dots \times m_d}$

Parameterized PDF :  $X_{i,\dots,id} \sim p(X_{i,\dots,id} | \theta_{i,\dots,id})$   
where  $\ell(\theta_{i,\dots,id}) = M_{i,\dots,id}$

$\ell : \mathbb{R} \rightarrow \mathbb{R}$  is an invertible link function connecting the model parameters  $M_{i,\dots,id}$  to the "natural" parameter of the distribution.

Model parameters :  $M \in \mathbb{R}^{m_1 \times \dots \times m_d}$

Goal : Find  $M$  that maximizes the likelihood of the observed entries of  $X$ .  
 $\hookrightarrow \Omega \subseteq [m_1] \times \dots \times [m_d]$

$$\max_M L(M; X, \Omega) = \prod_{(i,\dots,id) \in \Omega} p(X_{i,\dots,id} | \theta_{i,\dots,id})$$

$$\text{where } \ell(\theta_{i,\dots,id}) = M_{i,\dots,id}$$

→ We only have (at most) one observation for each entry of  $X$ !

We will assume a low rank structure of  $M$  (interdependence between the  $\theta_{i,\dots,id}$ )  
 $\hookrightarrow$  low CP rank of  $M$ .

Reformulation of optimization problem:

$$\min_M F(M; X, \Omega) = \sum_{(i,\dots,id) \in \Omega} f(X_{i,\dots,id}, M_{i,\dots,id}) \text{ s.t. } \text{rank}_{\text{CP}}(M) \leq R$$

$$\text{where } f(x, m) = -\log p(x | \ell^{-1}(m))$$

$\hookrightarrow$  "loss function" (assumed differentiable)

## Examples of data distribution

### ① Gaussian distribution

(Gaussian MLE  $\hat{\theta} \approx$  Mean squared error)

$$X_{i,\dots,d} = M_{i,\dots,d} + \varepsilon_{i,\dots,d} \text{ where } \varepsilon_{i,\dots,d} \sim \mathcal{N}(0, \sigma^2) \quad (\text{constant})$$

Equivalently:  $X_{i,\dots,d} \sim \mathcal{N}(\theta_{i,\dots,d}, \sigma^2)$  where  $\theta_{i,\dots,d} = M_{i,\dots,d}$   
the link function is the identity  $\ell(\theta) = \theta$

Simple derivation:

$$f(x, m) = -\log \mathcal{N}(x; m, \sigma^2) = \frac{(x-m)^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)$$

↳ ignoring  $\sigma^2$  and constants

$$\boxed{f(x, m) = (x-m)^2}$$

$\Rightarrow \star$  is then equivalent to CP decomposition.

### ② Bernoulli distribution

If  $X_{i,\dots,d} \in \{0, 1\}$  (binary data), we can assume  $X_{i,\dots,d} \sim \mathcal{B}(\theta_{i,\dots,d})$

$$p(x|\theta) = \theta^x (1-\theta)^{1-x}, x \in \{0, 1\}, \theta \in [0, 1]$$

link function?

• identity as we need to constrain  $m \in [0, 1]$

• odds ratio:  $m = \ell(\theta) = \frac{\theta}{1-\theta}$

↳ only need to enforce  $m \geq 0$

• log odds ratio:

$$m = \ell(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$$

$$\hookrightarrow \theta = \frac{e^m}{1+e^m}, 1-\theta = \frac{1}{1+e^m}$$

$$\dots f(x, m) = -\log p(x|\theta) = -x \log \theta - (1-x) \log(1-\theta)$$

cross-entropy

$$\boxed{f(x, m) = \log(1+e^m) - xm}$$

### ③ Other distributions

Table 1: Statistically-motivated loss functions. Parameters in blue are assumed to be constant. Numerical adjustments are indicated in red.

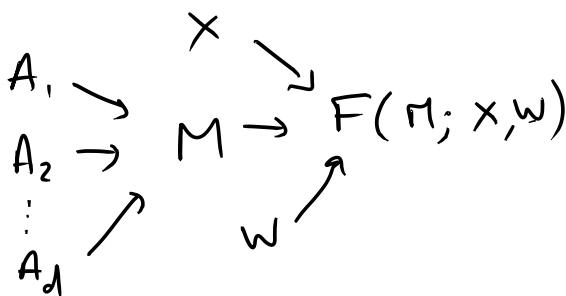
Distribution	Link function	Loss function	Constraints
$\mathcal{N}(\mu, \sigma)$	$m = \mu$	$(x-m)^2$	$x, m \in \mathbb{R}$
Gamma( $k, \sigma$ )	$m = k\sigma$	$x/(m+\epsilon) + \log(m+\epsilon)$	$x > 0, m \geq 0$
Rayleigh( $\theta$ )	$m = \sqrt{\pi/2}\theta$	$2\log(m+\epsilon) + (\pi/4)(x/(m+\epsilon))^2$	$x > 0, m \geq 0$
Poisson( $\lambda$ )	$m = \lambda$	$m - x \log(m+\epsilon)$	$x \in \mathbb{N}, m \geq 0$
	$m = \log \lambda$	$e^m - xm$	$x \in \mathbb{N}, m \in \mathbb{R}$
Bernoulli( $\rho$ )	$m = \rho / (1-\rho)$	$\log(m+1) - x \log(m+\epsilon)$	$x \in \{0, 1\}, m \geq 0$
	$m = \log(\rho / (1 - \rho))$	$\log(1+e^m) - xm$	$x \in \{0, 1\}, m \in \mathbb{R}$
NegBinom( $r, \rho$ )	$m = \rho / (1-\rho)$	$(r+x) \log(1+m) - x \log(m+\epsilon)$	$x \in \mathbb{N}, m \geq 0$

### Gradient based optimization

★ :  $\min_{A_1, \dots, A_d} F(M; X, W) = \sum_{i_1, \dots, i_d} w_{i_1, \dots, i_d} f(x_{i_1, \dots, i_d}, M_{i_1, \dots, i_d})$

s.t.  $M = [A_1, A_2, \dots, A_d]$

$\begin{matrix} m_1 \times R \\ m_2 \times R \\ \vdots \\ m_d \times R \end{matrix}$



We want to show

$$\frac{\partial F(M; X, W)}{\partial A_k} = -y^{(k)} \underbrace{\begin{pmatrix} A_1 & & & \\ & A_{k-1} & & \\ & & I & \\ & & & R \end{pmatrix}}_{R \times R} \quad \text{where } y \in \mathbb{R}^{m_1 \times \dots \times m_d}$$

$y_{i_1, \dots, i_d} = w_{i_1, \dots, i_d} \frac{\partial f(x_{i_1, \dots, i_d}, M_{i_1, \dots, i_d})}{\partial M_{i_1, \dots, i_d}}$

$1 \times 1$

of the same size as  $X$  but it is sparse if  $W$  is sparse.

$$= y^{(k)} (A_d \odot \dots \odot A_{k+1} \odot A_{k-1} \odot \dots \odot A_1)$$

$$\frac{\partial F(M; X, w)}{\partial A_k} = \underbrace{\frac{\partial F(M; X, n)}{\partial M}}_{1 \times m_k \times R} \frac{\partial M}{\partial A_k}$$

$$1 \times m_k \times R \quad 1 \times m_1 \times \dots \times m_d \quad m_1 \times \dots \times m_d \times m_k \times R$$

$$\begin{aligned} \frac{\partial F(M; X, n)}{\partial M}_{i, \dots, id} &= \frac{\partial \left( \sum_{j_1 \dots j_d} w_{j_1 \dots j_d} f(x_{j_1 \dots j_d}, M_{j_1 \dots j_d}) \right)}{\partial M_{i, \dots, id}} \\ &= w_{i, \dots, id} \frac{\partial f(x_{i, \dots, id}, M_{i, \dots, id})}{\partial M_{i, \dots, id}} := y_{i, \dots, id} \end{aligned}$$

$$\frac{\partial M}{\partial A_k} = \frac{\partial \begin{pmatrix} -A_1 & & \\ -A_2 & \curvearrowright I & \\ \vdots & & \\ -A_d & & \end{pmatrix}}{\partial A_k} = \begin{array}{c} m_1 \text{---} A_1 \\ \vdots \\ m_{k-1} \text{---} A_{k-1} \\ m_k \text{---} \color{red} R \\ \vdots \\ m_{k+1} \text{---} A_{k+1} \\ \vdots \\ m_d \text{---} A_d \end{array}$$

$m_1 \times \dots \times m_d \times m_k \times R$

$$\frac{\partial F(M; X, w)}{\partial A_k} = \frac{\partial F(M; X, n)}{\partial M} \frac{\partial M}{\partial A_k}$$

$$= \frac{-y}{m_k} \begin{array}{c} A_1 \\ \vdots \\ A_{k-1} \\ I \\ A_{k+1} \\ \vdots \\ A_d \end{array} = y_{(k)} (A_d \odot \dots \odot A_{k+1} \odot A_{k-1} \odot \dots \odot A_1)$$

Follow up paper:

Stochastic Gradients for Large-Scale Tensor Decomposition\*

Tamara G. Kolda<sup>†</sup> and David Hong<sup>‡</sup>